



INVESTOR IN PEOPLE

The Patent Office  
Concept House  
Cardiff Road

Newport REC'D 19 OCT 2004  
South Wales  
NP10 8QW/PO PCT

I, the undersigned, being an officer duly authorised in accordance with Section 74(1) and (4) of the Deregulation & Contracting Out Act 1994, to sign and issue certificates on behalf of the Comptroller-General, hereby certify that annexed hereto is a true copy of the documents as originally filed in connection with the patent application identified therein.

In accordance with the Patents (Companies Re-registration) Rules 1982, if a company named in this certificate and any accompanying documents has re-registered under the Companies Act 1980 with the same name as that with which it was registered immediately before re-registration save for the substitution as, or inclusion as, the last part of the name of the words "public limited company" or their equivalents in Welsh, references to the name of the company in this certificate and any accompanying documents shall be treated as references to the name with which it is so re-registered.

In accordance with the rules, the words "public limited company" may be replaced by p.l.c., plc, P.L.C. or PLC.

Re-registration under the Companies Act does not constitute a new legal entity but merely subjects the company to certain additional company law rules.

BEST AVAILABLE COPY



Signed

Dated 12 October 2004

**PRIORITY  
DOCUMENT**

SUBMITTED OR TRANSMITTED IN  
COMPLIANCE WITH RULE 17.1(a) OR (b)

Patents Act 1977  
(Rule 16)

30 SEP 2003

LONDON

010CT03 E841126-1 D03052  
P01/7700 0.00-0322899.6

The Patent Office  
Cardiff Road  
Newport  
Gwent NP10 8QQ

**Request for grant of a patent**

*(See the notes on the back of this form. You can also get an explanatory leaflet from the Patent Office to help you fill in this form)*

1. Your reference

**A30425**

2. Patent application number  
*(The Patent Office will fill in this part)*

30 SEP 2003

**0322899.6**

3. Full name, address and postcode of the or of each applicant *(underline all surnames)*

**BRITISH TELECOMMUNICATIONS public limited company  
81 NEWGATE STREET  
LONDON, EC1A 7AJ, England  
Registered in England: 1800000**

Patents ADP number *(if you know it)*

**1867002**

If the applicant is a corporate body, give the country/state of its incorporation

**UNITED KINGDOM**

4. Title of the invention

**INFORMATION RETRIEVAL**

5. Name of your agent *(if you have one)*

"Address for Service" in the United Kingdom to which all correspondence should be sent *(including the postcode)*

**BT GROUP LEGAL  
INTELLECTUAL PROPERTY DEPARTMENT  
HOLBORN CENTRE  
120 HOLBORN  
LONDON, EC1N 2TE**

859 191 9001

Patents ADP number *(if you know it)*

**4867001**

6. If you are declaring priority from one or more earlier patent applications, give the country and the date of filing of the or of each of these earlier applications and *(if you know it)* the or each application number

Country

Priority application number  
*(if you know it)*

Date of filing  
*(day / month / year)*

7. If this application is divided or otherwise derived from an earlier UK application, give the number and the filing date of the earlier application

Number of earlier application

Date of filing  
*(day/month/year)*

8. Is a statement of inventorship and of right to grant of a patent required in support of this request? *(Answer 'Yes' if:*

**YES**

- a) any applicant named in part 3 is not an inventor, or
- b) there is an inventor who is not named as an applicant, or
- c) any named applicant is a corporate body.

*(See note (d))*

9. Enter the number of sheets for any of the following items you are filing with this form. Do not count copies of the same document

Continuation sheets of this form -

Description - 11

Claim(s) - 1

Abstract - 1

Drawing(s) - 2 + 2 SN

10. If you are also filing any of the following, state how many against each item

Priority Documents

Translations of priority documents

Statement of inventorship and right to grant of a patent (Patents Form 7/77)

Request for preliminary examination and search (Patents Form 9/77)

Request for substantive examination (Patents Form 10/77)

Any other documents (please specify)

11. I/We request the grant of a patent on the basis of this application.

Signature(s)

Date:

30 September 2003

LLOYD, Barry George William, Authorised Signatory

12. Name and daytime telephone number of person to contact in the United Kingdom

Rohini R RANJITKUMAR

020 7492 8456

### Warning

*After an application for a patent has been filed, the Comptroller of the Patent Office will consider whether publication or communication of the invention should be prohibited or restricted under Section 22 of the Patents Act 1977. You will be informed if it is necessary to prohibit or restrict your invention in this way. Furthermore, if you live in the United Kingdom, Section 23 of the Patents Act 1977 stops you from applying for a patent abroad without first getting written permission from the Patent Office unless an application has been filed at least 6 weeks beforehand in the United Kingdom for a patent for the same invention and either no direction prohibiting publication or communication has been given, or any such direction has been revoked.*

### Notes

- If you need help to fill in this form or you have any questions, please contact the Patent Office on 0645 500505.
- Write your answers in capital letters using black ink or you may type them.
- If there is not enough space for all the relevant details on any part of this form, please continue on a separate sheet of paper and write "see continuation sheet" in the relevant part(s). Any continuation sheet should be attached to this form.
- If you have answered 'Yes' Patents Form 7/77 will need to be filed.
- Once you have filled in the form you must remember to sign and date it.
- For details of the fee and ways to pay please contact the Patent Office.

## INFORMATION RETRIEVAL

This invention relates to information retrieval and in particular to a method and apparatus for determining similarity of words and information content of documents as an aid to information retrieval.

There are a number of known techniques by which semantic similarity of documents may be determined. In one such technique, a document is represented by a vector, each value in the vector being a measure of the incidence of a corresponding word or term in the document. A measure of semantic similarity between two such documents may then be calculated as the scalar product, also known as the dot product, of the corresponding document vectors. Such a measure of document similarity forms the basis of a known document clustering technique whereby documents having semantically similar content may be assembled into groups of documents apparently relating to similar subject matter. However, by this technique, the measure of semantic similarity between two documents is based only upon those words or terms that occur in both documents. That is, document vectors must relate to the same set of words or terms. One problem with this technique is that when two documents describe the same topic but use slightly different terminology, the technique would fail to recognise the semantic similarity.

According to a first aspect of the present invention there is provided a method of generating a word similarity matrix for a plurality of words, selected from a set of documents, for use in determining semantic similarity of documents in said set, comprising the steps of:

- (i) for each word of said plurality of words:
    - (a) identifying, in documents of said set, distinct word sequences comprising the word and a predetermined number of other words;
    - (b) calculating a relative frequency of occurrence for each distinct word sequence among word sequences containing the word; and
    - (c) generating a fuzzy set comprising, for groups of word sequences containing the word, corresponding fuzzy membership values calculated from the relative frequencies determined at step (b);
  - (ii) calculating, for each pair of words of said plurality of words, using respective fuzzy sets generated at step (i), a probability that the first word of the pair is semantically suitable as a replacement for the second word of the pair; and
  - (iii) creating a word similarity matrix comprising said calculated probabilities.
- Preferably, the method comprises the further step of:

(iv) adding a new document to said set of documents and, using a set of words selected from said new document, performing an incremental update to said word similarity matrix by means of steps (i) and (ii) performed in respect of said selected words using corresponding word sequences identified in said new document.

5 According to a second aspect of the present invention, there is provided an information retrieval apparatus for use in retrieving information from a set of documents, comprising:

an input for receiving a search query;

10 a word similarity matrix generated according to the first aspect of the present invention;

query enhancement means for modifying terms of a received search query with reference to said word similarity matrix; and

information retrieval means for searching said set of documents for relevant information using a query enhanced by said query enhancement means.

15 Preferred embodiments of the present invention will now be described in more detail and with reference to the accompanying drawings, of which:

Figure 1 is an overview of a process according to preferred embodiments of the present invention;

20 Figure 2 is a flow chart showing steps in a process for generating a word replaceability matrix according to a preferred embodiment of the present invention.

An overview of a preferred process according to embodiments of the present invention will firstly be described with reference to Figure 1.

Referring to Figure 1, a diagram is shown representing, in overview, a process for analysing a document set 100 in order to generate a word replaceability matrix 115. The contents of the document set 100 are analysed 105 to calculate a measure of the semantic similarity of words used in the document set. The determined similarities are checked in a verification process 110 and the calculated values of the measure are stored as a word replaceability matrix 115, each value being indicative of the degree of semantic similarity between a respective pair of words and hence the probability that the first word of the pair is suitable for use in place of the second word. The matrix 115 may be exploited in a number of ways of which two are shown in Figure 1: to cluster 120 documents in the document set 100 into distinct information categories; and enhancement 125 of search queries for use in information retrieval in the document set 100 or in other document sets. Both applications will be discussed in more detail below.

In a preferred embodiment of the present invention to be described below, a particular technique is used to calculate the semantic similarity of words used in the input document set 100. The technique is based upon identification of so-called "n-grams" of words occurring in the document set 100, where an n-gram is any sequence of n consecutive words occurring in a document. For example the sequence of words "the cat is blue" is a 4-gram of words. The main purpose of identifying n-grams in the present invention is to understand and to represent the context in which particular words are used in a document. The value of n is determined at the outset, and the inventors in the present case have found that a value  $n=3$  or  $n=4$  gives good results, although other values may also be selected. However, use of significantly higher values of n does not appear to improve the performance of the technique. For each word, a fuzzy set of corresponding n-grams is formed based upon the observed probabilities of each n-gram occurring in the document set 100. The technique of semantic unification, described for example in a paper by J. F. Baldwin, J. Lawry, and T. P. Martin: "Efficient Algorithms for Semantic Unification", in *Proc. Information Processing and the Management of Uncertainty*, 1996, Spain, is then used calculate the semantic similarity of words from their respective fuzzy sets and hence to determine the probability that one word would be a suitable replacement for another. The calculated probabilities and the respective word pairs are collated into a table to form a so-called "Word Replaceability Matrix" 115. Preferably, a predetermined threshold is applied so that only those probabilities that exceed the threshold, and hence only the strongest respective word similarities, are recorded in the matrix 115.

A preferred process 105 for calculating the semantic similarity of words occurring in the document set 100 and hence for generating a word replaceability matrix 115 will now be described in detail with reference to Figure 2, according to a preferred embodiment of the present invention.

Referring to Figure 2, the process begins, and at STEP 200 a document set 100 is input, the set comprising a number of documents containing readable text, for example documents in ASCII plain text format, or XML files. At STEP 205 some initial, optional, word analysis may be carried out by way of an initial filtering step to eliminate certain types of word from consideration in the remaining steps in this process and hence to select a first set of words as candidates for representation in a resultant Word Replaceability Matrix 115. STEP 205 is considered to be optional because it is not essential to the working of the present invention to limit the choice of words represented in the resultant word replaceability matrix 115. However, there are certain advantages to

eliminating certain types of word from the remaining steps in this process, not least in savings in processing required to generate the matrix 115. Certain types of "low value" word are unlikely to be useful in a resultant matrix 115, words such as "a", "the", "and", "or". In addition, the inventors have found that there is little advantage to including other words that occur either very frequently or very infrequently in the input document set 100. Thus, processing at STEP 205 may include an analysis of the frequency of occurrence of each word in the input document set 100 and the elimination from further consideration of those words having a frequency of occurrence lying in the first and fourth quartiles of the observed frequency distribution. However, this latter step may be omitted in particular when carrying out incremental updates to the matrix 115, triggered by the input of a further document for example.

At STEP 210 those words remaining after STEP 205 may be processed in a word stemming algorithm, a suitable stemming algorithm being the Porter Stemmer algorithm, as described in M. F. Porter: An Algorithm for Suffix Stripping, Automated Library and Information Systems, Vol. 14, No. 3, pp. 130-137, 1980.

At STEP 215, for each word output from STEP 210, the input documents 100 are analysed to identify the n-grams of surrounding words, each n-gram being representative of a context in which the word is being used. The value of n is predetermined and a value of 3 or 4 has been found by the inventors in the present case to give satisfactory results. Preferably, in identifying the n-grams, characters such as punctuation marks, brackets, inverted commas, hyphens and underscores are ignored, and n-grams are not selected where they overlap sentence boundaries. Formally, the following natural language procedure is followed to identify n-grams in a document:

```

25      DOC =      START      WORDS?      END
      WORDS =  WORD | WORD SPACE* WORD
      WORD =    (any char not {' '}) *
      START =   (start of file)
      END =     (end of file)
30      SPACE =  white space or ,.;:-
      Ignore characters  ""(){}[]
      Also ignore n-grams that would contain a "." at a position other than at the end

```

Consider, by way of example, the following four sentences, found in an input document set to contain the word *brown*:

- The quick brown fox jumps over the lazy dog.
  - The quick brown cat jumps onto the active dog.
  - The slow brown fox jumps onto the quick brown cat.
- 5 - The quick brown cat leaps over the quick brown fox.

Assuming that a value of  $n=3$  has been chosen for this operation of the process, then the word *brown* occurs in three distinct contexts represented by the 3-grams formally denoted by

10

brown: (quick,fox)

brown: (quick,cat)

brown: (slow,fox)

15

At STEP 220, for each word, the relative frequency of occurrence of each corresponding n-gram is calculated. That is, for each word, the frequency of occurrence of each distinct and corresponding n-gram is divided by the total number of n-grams containing the word to give, for each distinct n-gram, a measure of the probability that the word appears in the document set 100 in the context represented by that distinct and

20

corresponding n-gram. To illustrate this, continuing with the example from STEP 215, the word *brown* occurs in a total of six 3-grams, represented by three distinct 3-grams having a frequency of occurrence shown in the following table:

brown.	total = 6		
	quick	fox	2
	quick	cat	3
	slow	fox	1

- 25 From this, the respective probabilities can be calculated to give the following probability distribution for the contexts of *brown*, in order of decreasing probability:

$$\text{Pr} \{ (\text{quick}, \text{cat}) \} = 1/2$$

$$\text{Pr} \{ (\text{quick}, \text{fox}) \} = 1/3$$

30

$$\text{Pr} \{ (\text{slow}, \text{fox}) \} = 1/6$$



At STEP 225, the probability values calculated at STEP 220 are used to generate a fuzzy set for each word. That is, for each distinct n-gram, or context of a word, the corresponding probability values are used to calculate fuzzy membership values for the word. Preferably, in calculating these fuzzy membership values, the underlying principle of "least prejudiced distribution" of probability mass is applied, meaning that in the absence of any bias towards one or other element in a group of n-grams, the probability mass assigned to the group is distributed equally amongst the composite n-grams. The principles of fuzzy membership values and probability mass assignment are described for example in J. F. Baldwin (1992) in "The Management of Fuzzy and Probabilistic Uncertainties for Knowledge-based Systems.", the *Encyclopedia of AI*, edited by S. A. Shapiro, published by John Wiley (2<sup>nd</sup> edition), pages 528-537.

This step in the process may be illustrated by a continuation of the example from STEP 220. Starting with the probabilities calculated at STEP 220 of the word *brown* arising in the document set in each of the 3-gram contexts as follows:

$$\text{Pr} \{ (\text{quick}, \text{cat}) \} = 1/2$$

$$\text{Pr} \{ (\text{quick}, \text{fox}) \} = 1/3$$

$$\text{Pr} \{ (\text{slow}, \text{fox}) \} = 1/6$$

and representing the corresponding fuzzy membership values to be determined as x, y and z respectively, then the assignment of probability mass across the possible contexts for the word *brown* would be represented by

$$\{(\text{quick}, \text{cat})\}: x-y, \{(\text{quick}, \text{cat}), (\text{quick fox})\}: y-z, \{(\text{quick}, \text{cat}), (\text{quick fox}), (\text{slow}, \text{fox})\}: z$$

In the absence of any bias in favour of one context over another, the probability masses y-z and z are assumed to be distributed evenly over the contexts in their respective groups. This distribution is therefore referred to as the least prejudiced distribution of the probability mass. While other distributions of the probability masses y-z and z are possible in general, no other distributions are considered in the present patent application.

On the assumption of a least prejudiced distribution of the probability mass, the fuzzy membership values for each context would therefore be required to satisfy the following equations, relating the fuzzy membership values to the calculated probabilities above:

$$(\text{quick}, \text{fox}): x - y + (y - z)/2 + z/3 = 1/2$$

$$(\text{quick}, \text{cat}): (y - z)/2 + z/3 = 1/3$$

$$(\text{slow}, \text{fox}): z/3 = 1/6$$

5 Solving these three simultaneous equations for  $x$ ,  $y$  and  $z$  gives fuzzy membership values of  $x=1$ ,  $y=5/6$  and  $z=1/2$ . Therefore the fuzzy set for the word *brown* is

$$\{(\text{quick}, \text{cat}) : 1, (\text{quick}, \text{fox}) : 0.833, (\text{slow}, \text{fox}) : 0.5\}$$

10 By this technique, fuzzy sets are generated for each of the words output from STEP 210 for which contexts (n-grams) were identified at STEP 215.

At STEP 230, for each pair of words, the corresponding fuzzy sets are used to calculate the probability that one word of the pair may be a semantically suitable word to use in place of the other word of the pair. These probabilities will ultimately be the basis of  
15 the Word Replaceability Matrix 115. The technique of point semantic unification is applied to calculate these probabilities from the membership values in the respective word fuzzy sets. However, to illustrate the principle, the example will be continued from STEP 225.

For the word *brown*, the following fuzzy set was generated at STEP 225:

20  $\{(\text{quick}, \text{cat}) : 1, (\text{quick}, \text{fox}) : 0.833, (\text{slow}, \text{fox}) : 0.5\}$

The mass assignment for the word *brown* is therefore

$$m(\text{brown}) = \{(\text{quick}, \text{cat})\}: 1/6, \{(\text{quick}, \text{cat}), (\text{quick}, \text{fox})\}: 1/3, \{(\text{quick}, \text{cat}), (\text{quick}, \text{fox}), (\text{slow}, \text{fox})\}: 1/2$$

25

Suppose that for another word, *black*, the following fuzzy set was generated at STEP 225:

$$\{(\text{quick}, \text{cat}) : 1, (\text{slow}, \text{fox}) : 0.75\}$$

30

The mass assignment for the word *black* is therefore

$$m(\text{black}) = \{(\text{quick}, \text{cat})\}: 1/4, \{(\text{quick}, \text{cat}), (\text{slow}, \text{fox})\}: 3/4$$

The degree of support for the word *black* being a semantically suitable replacement given the word *brown* may be represented in table form as follows, where the mass assignments for the given word *brown* are arranged across the columns of the table and those for the potential replacement word *black* being arranged as the rows:

5

	<b>{{(quick,cat)}}: 1/6</b>	<b>{{(quick,cat), (quick,fox)}}: 1/3</b>	<b>{{(quick,cat), (quick,fox), (slow,fox)}}: 1/2</b>
<b>{{(quick,cat)}}: 1/4</b>	$1/4 \times 1/6$	$1/2 \times 1/4 \times 1/3$	$1/3 \times 1/4 \times 1/2$
<b>{{(quick,cat), (slow,fox)}}: 3/4</b>	$3/4 \times 1/6$	$1/2 \times 3/4 \times 1/3$	$2/3 \times 3/4 \times 1/2$

The probability of the word *black* being a suitable replacement for the word *brown* is the sum of the values in the table, giving a conditional probability  $\Pr(\text{black} \mid \text{brown}) =$   
 10 0.625. Similarly, the probability of the word *brown* being a suitable replacement for the word *black* may be calculated using a corresponding table as the conditional probability  $\Pr(\text{brown} \mid \text{black}) = 0.8125$ .

By performing these calculations for each pair of words for which fuzzy sets were generated at STEP 225, a table of conditional probabilities is generated. Preferably, a  
 15 predetermined threshold is applied so that only those conditional probabilities that exceed the threshold, and hence only the strongest respective word similarities, are preserved in the table, all other probabilities being set to zero.

At STEP 235, a verification step (110) may be performed to automatically or semi-automatically eliminate any of the more unlikely relationships identified between  
 20 words under this process 105. In a preferred method, a lexical database such as Wordnet™, accessible over the Internet at <http://www.cogsci.princeton.edu/~wn/>, may be used in a procedure to check the semantic relationships identified and, if necessary, to modify corresponding probability values in the table generated at STEP 225, setting them to zero for example where a relationship is apparently invalid. For example, a process  
 25 may be executed whereby each word in the table is submitted in turn to Wordnet and a corresponding list of synonyms, hyponyms, hypernyms and antonyms is returned. For each word in the generated table having a calculated conditional probability in excess of a predetermined threshold, a comparison is made with the semantic relationship suggested by the list returned by Wordnet. If there is no apparent semantic relationship suggested by

Wordnet, or if the meanings of the words are clearly opposite, then the replaceability suggested by the calculated value of conditional probability in the table is likely to be false and the value may be overwritten with a zero. Where the result of the comparison is not clear-cut, a manual verification may be carried out, achieved preferably by presenting to a user, as background to the apparent relationship between the words, the respectively generated fuzzy sets.

The table resulting from verification STEP 235 (110) is the Word Replaceability Matrix 115.

Once the matrix 115 has been generated it may be exploited in a number of ways. For example, the word replaceability matrix 115 may be used in an enhancement to the known vector dot product technique for assessing semantic similarity of documents, described above in the introductory part of the present patent specification. A weakness of that known vector dot product technique is that related documents that use different terminology are not identified as being semantically related. The enhancement made possible by the word replaceability matrix 115 of the present invention allows the measure of similarity to be based upon words that are not necessarily the same between documents but which are nevertheless semantically similar to some degree.

In the known vector dot product technique, if a first document is represented by a document vector  $\underline{v}_1 = (v_{11}, v_{12}, \dots, v_{1k})$  and a second document is represented by a document vector  $\underline{v}_2 = (v_{21}, v_{22}, \dots, v_{2k})$ , where the values  $v_{ij}$  are indicative of the incidence of a j-th word of a common set of k words in the document i, then the dot product

$$\underline{v}_1 \cdot \underline{v}_2 = \sum_i v_{1i} v_{2i}$$

provides a value indicative of the semantic similarity between the documents. However, by using the probability values in the word replaceability matrix 115, this known measure of semantic similarity may be enhanced so that not only are identical words considered in the calculation of document similarity, but also other words represented in the word replaceability matrix 115 that may be semantically related.

Assuming that the word replaceability matrix 115 contains m words, so that the matrix 115 is an m x m matrix of probability values, then for an i-th document, a 1 x m matrix of values  $u_{ij}$  may be formed where the j-th value is indicative of the frequency of occurrence, in the i-th document, of the j-th word in the matrix 115. If a particular word of the matrix does not occur in the document, then a value of zero appears in the corresponding position in the 1 x m matrix for that document. The values  $u_{ij}$  in the 1 x m

matrix are normalised so that a document containing an unusually high proportion of the words represented in the matrix does not skew the calculation that follows.

The semantic similarity  $S_{12}$  between a first document, represented by a  $1 \times m$  matrix

5

$$\underline{u}_1 = (u_{11} \ u_{12} \ \dots \ u_{1m})$$

and a second document, represented by a  $1 \times m$  matrix

10

$$\underline{u}_2 = (u_{21} \ u_{22} \ \dots \ u_{2m}),$$

is calculated, by this enhanced measure, according to the following multiplication of matrices

15

$$S_{12} = \sum_j \sum_i w_{ji} u_{1i} u_{2j}$$

where  $w_{ji}$  is the probability, read from the matrix 115, that the  $j$ -th word represented in the matrix 115 is semantically suitable as a replacement for the  $i$ -th word of the matrix 115.

Using this enhanced measure of semantic similarity between documents, for example those in document set 100, documents may be clustered (120) into groups of documents having related information content.

Of course, the matrix 115 may be used as a semantic dictionary both in relation to documents of the document set 100 on which it was based, or in relation to other documents. However, a particular advantage of the process 105, 110 described above with reference to Figure 2 is that the matrix 115 may be incrementally updated as new documents, and hence new words, are considered. In particular, on adding a new document, processing steps 205 to 235 of Figure 2, as described above, may be performed on the basis of a set of words, optionally selected from the new document at STEP 205, by generating fuzzy sets for words not already represented in the matrix 115 and by updating the fuzzy sets for those words of the new document that are represented in the matrix 115. The fuzzy sets for the new words are generated at STEP 225 entirely on the basis of  $n$ -grams identified at STEP 215 in the new document. The fuzzy membership values in the fuzzy sets for those selected words already represented in the matrix 115 are updated, at STEP 225, having included any new distinct  $n$ -grams identified at STEP

215 from the new document and having updated the probabilities, at STEP 220, both for the existing n-grams and for the new n-grams. Corresponding entries in the matrix 115 are then recalculated at STEP 230 in respect of the updated words and the matrix is extended as necessary with any new words selected from the new document.

5           As mentioned above with reference to Figure 1, besides application to an improved document clustering technique (120), the word replaceability matrix 115 may be used to extend or modify terms in a user's search query for use in an information retrieval system. In particular, a set of words entered by a user may be extended with semantically similar words identified with reference to the matrix 115 in order to improve the chances of  
10 a search engine returning a more complete set of relevant documents. This is likely to be particularly effective when searching for information contained in the document set on which the matrix 115 was based, although as more documents are considered and as the matrix 115 is incrementally updated, the more broadly-based semantic relationships and the increased number of words represented in the matrix 115 make it increasingly useful  
15 as a semantic dictionary for improving the information retrieval performance of search engines with respect to other information sets.

## CLAIMS

1. A method of generating a word similarity matrix for a plurality of words, selected from a set of documents, for use in determining semantic similarity of documents in said set, comprising the steps of:
  - (i) for each word of said plurality of words:
    - (a) identifying, in documents of said set, distinct word sequences comprising the word and a predetermined number of other words;
    - (b) calculating a relative frequency of occurrence for each distinct word sequence among word sequences containing the word; and
    - (c) generating a fuzzy set comprising, for groups of word sequences containing the word, corresponding fuzzy membership values calculated from the relative frequencies determined at step (b);
  - (ii) calculating, for each pair of words of said plurality of words, using respective fuzzy sets generated at step (i), a probability that the first word of the pair is semantically suitable as a replacement for the second word of the pair; and
  - (iii) creating a word similarity matrix comprising said calculated probabilities.
2. A method according to Claim 1, further comprising the step of:
  - (iv) adding a new document to said set of documents and, using a set of words selected from said new document, performing an incremental update to said word similarity matrix by means of steps (i) and (ii) performed in respect of said selected words using corresponding word sequences identified in said new document.
3. An information retrieval apparatus for use in retrieving information from a set of documents, comprising:
  - an input for receiving a search query;
  - a word similarity matrix generated in respect of said set of documents by the method of Claim 1;
  - query enhancement means for modifying terms of a received search query with reference to said word similarity matrix; and
  - information retrieval means for searching said set of documents for relevant information using a query enhanced by said query enhancement means.

## ABSTRACT

## INFORMATION RETRIEVAL

5        A method and apparatus are provided for generating, from an input set of documents, a word replaceability matrix defining semantic similarity between words occurring in the input document set. For each word, distinct word sequences of predetermined length are identified from the documents of the set, each word sequence being indicative of the context in which the word was used and, according to the relative  
10 frequency of occurrence of the identified word sequences for the word, fuzzy sets are generated for each word comprising membership values for corresponding groups of word sequences. For each pair of words occurring in the document set, their respective fuzzy sets are used to calculate the probability that the first word of a pair is semantically suitable as a replacement for the second word of the pair, these probabilities being  
15 collated to form a word similarity matrix for use in an improved method of determining document similarity and in information retrieval.

Figure (2)



1/2

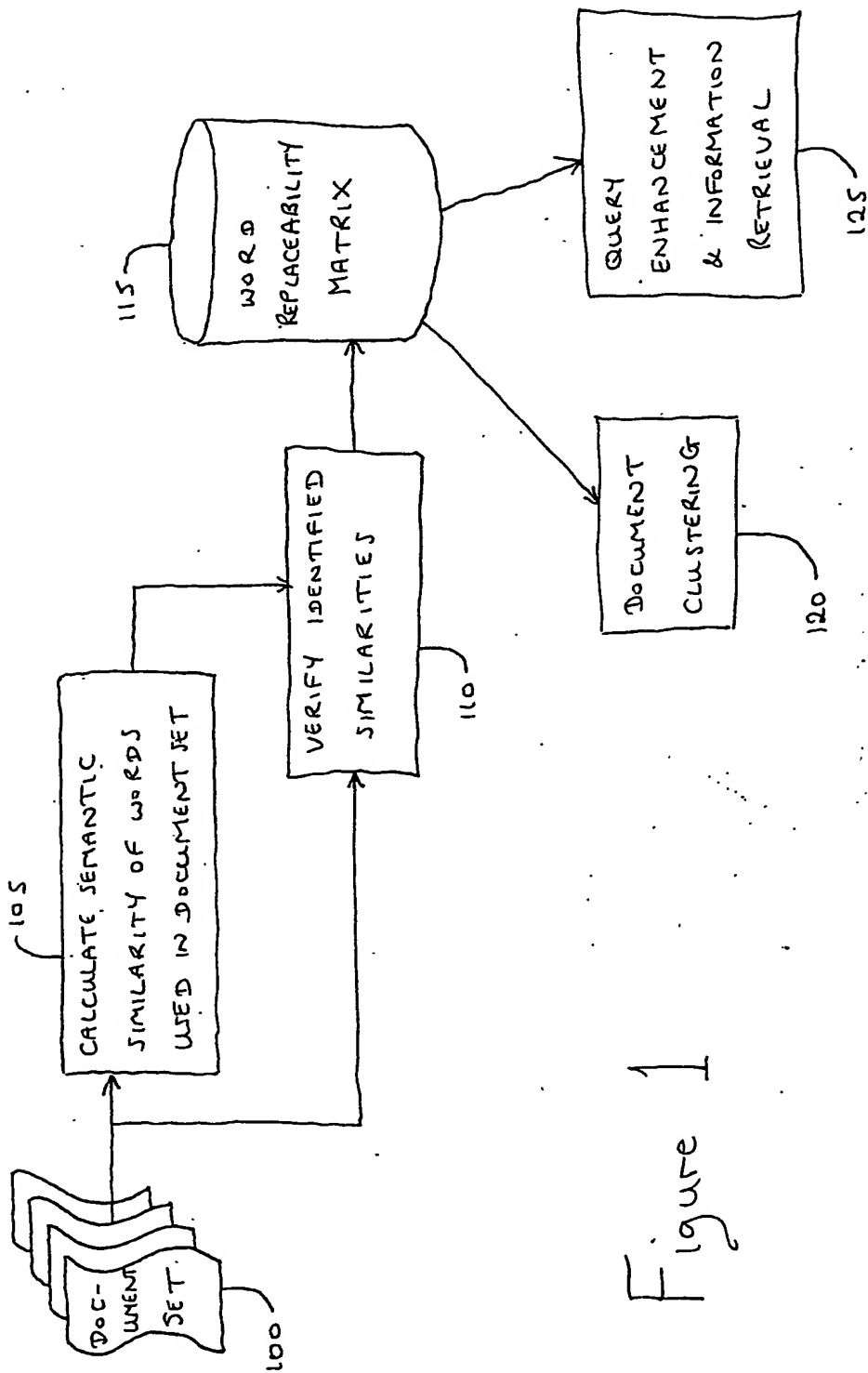


Figure 1

2/2

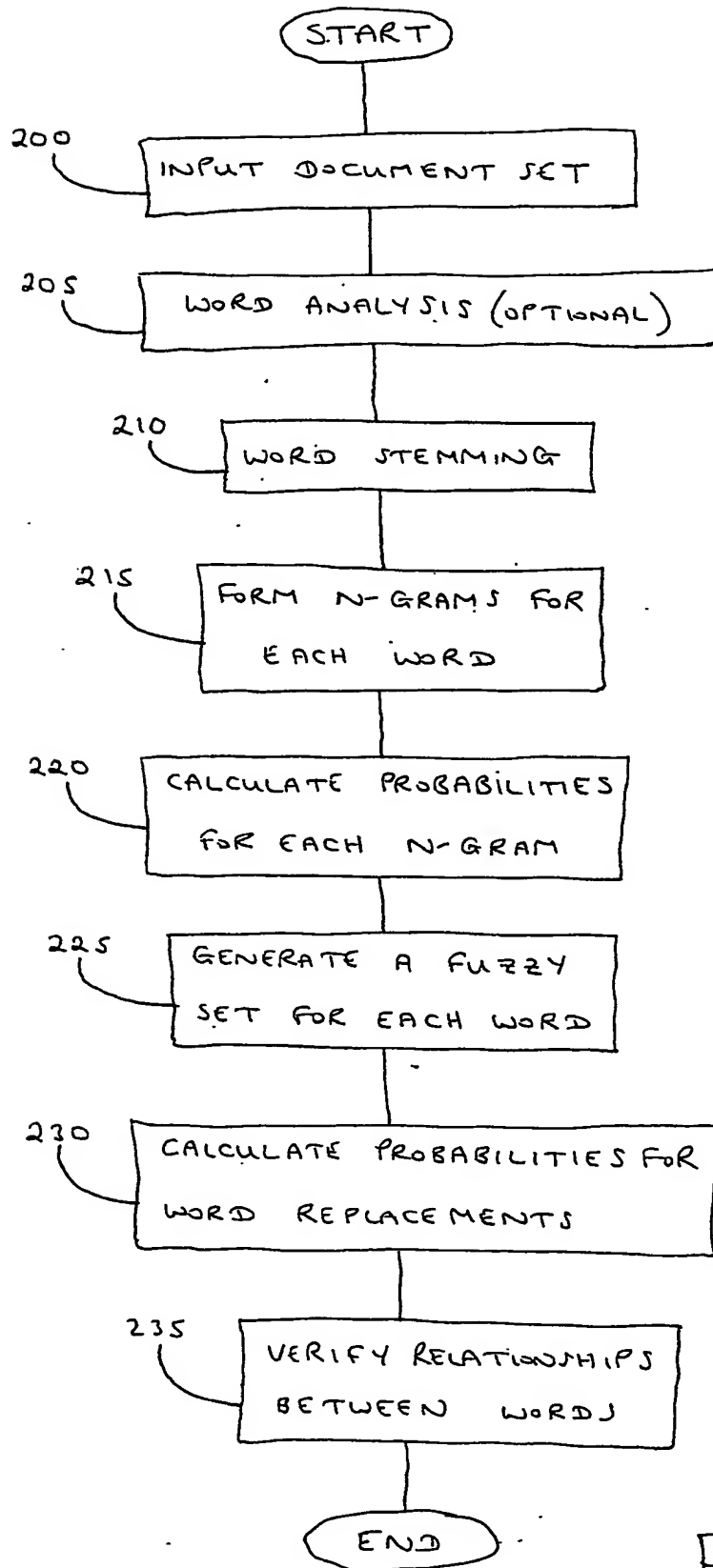


Figure 2

**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☒ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☒ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☒ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**